

# Methodology

## Lesson 4: Evaluating the models

---

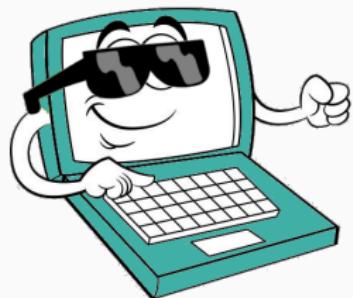
Carlos Ramisch

`first.last@lis-lab.fr`

M2 IAAA - based on the course *Zen Research*  
By Carlos Ramisch and Manon Scholivet

# The Awesome system

My brand new **Awesome** system...



should soon replace



...the old **Baseline** system!

# Expectation...

	dataset	metric1	metric2	metric3 <sup>1</sup>
Baseline system	DS1	82.3	75.9	48.0
Awesome system	DS1	<b>95.3</b>	<b>89.8</b>	<b>65.4</b>
Baseline system	DS2	67.7	65.2	56.8
Aweseome system	DS2	<b>80.3</b>	<b>91.1</b>	<b>69.8</b>
Baseline system	DS3	77.6	74.1	92.8
Awesome system	DS3	<b>84.9</b>	<b>78.3</b>	<b>98.1</b>

⇒ The Awesome system is **better** than the Baseline! 

<sup>1</sup>Higher is better

# ... Vs. reality!

	dataset	metric1	metric2	metric3
Baseline system	DS1	<b>82.3</b>	75.9	48.0
Awesome system	DS1	80.7	<b>76.2</b>	<b>50.4</b>
Baseline system	DS2	67.7	<b>65.2</b>	<b>56.8</b>
Awesome system	DS2	<b>67.9</b>	nan	49.6
Baseline system	DS3	77.6	<b>74.1</b>	92.8
Awesome system	DS3	<b>79.0</b>	<b>74.1</b>	<b>93.4</b>

⇒ Wake up and smell the coffee 😴

## Results analysis

- Identify overall trends
- Identify potential sources of problems (or bugs)
- Ensure conclusions are valid, claims are (statistically) sound

# Experimental results

- Diversity of experiments  $\implies$  diversity of results
  - Task at hand
  - Datasets
  - Evaluation metrics
  - ...
- This course: no silver bullet, rather a toolbox



# Outline

---

Statistics in a nutshell

Evaluation metrics

Statistical significance

Discussion

# Statistics

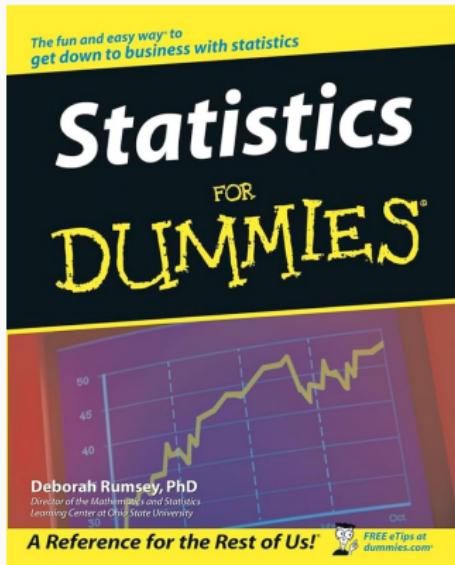
- A mathematical framework to analyse data
  - Foundations: probability theory
- Statistical inference  $\implies$  data science, machine learning
  - Also: finances, health, biology, physics, social sciences, ...
- Identify trends, test hypotheses, measure correlations, ...



# The problem with statistics

Finding good learning materials in statistics is hard

Too applied:



Too theoretical:

## Weak Law of Large Numbers

The weak law of large numbers (cf. the [strong law of large numbers](#)) is a result in probability theory also known as Bernoulli's theorem. Let  $X_1, \dots, X_n$  be a sequence of independent and identically distributed random variables, each having a [mean](#)  $\langle X_i \rangle = \mu$  and [standard deviation](#)  $\sigma$ . Define a new variable

$$X = \frac{X_1 + \dots + X_n}{n}.$$

Then, as  $n \rightarrow \infty$ , the sample mean ( $\bar{x}$ ) equals the population [mean](#)  $\mu$  of each variable.

$$\begin{aligned}\langle X \rangle &= \left\langle \frac{X_1 + \dots + X_n}{n} \right\rangle \\ &= \frac{1}{n} (\langle X_1 \rangle + \dots + \langle X_n \rangle) \\ &= \frac{n \mu}{n} \\ &= \mu.\end{aligned}$$

In addition,

$$\begin{aligned}\text{var}(X) &= \text{var}\left(\frac{X_1 + \dots + X_n}{n}\right) \\ &= \text{var}\left(\frac{X_1}{n}\right) + \dots + \text{var}\left(\frac{X_n}{n}\right) \\ &= \frac{\sigma^2}{n^2} + \dots + \frac{\sigma^2}{n^2} \\ &= \frac{\sigma^2}{n}.\end{aligned}$$

Therefore, by the [Chebyshev inequality](#), for all  $\epsilon > 0$ ,

$$P(|\bar{X} - \mu| \geq \epsilon) \leq \frac{\text{var}(\bar{X})}{\epsilon^2} = \frac{\sigma^2}{n \epsilon^2}.$$

# What usually happens

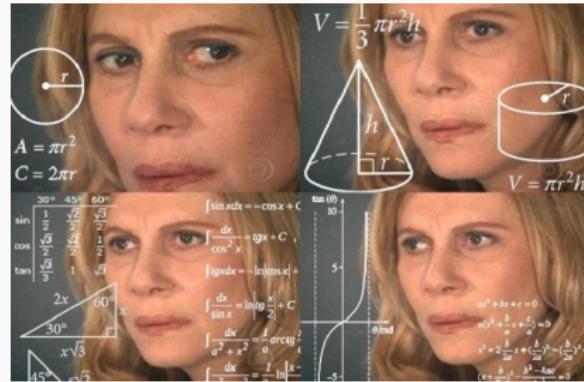
- A given statistical **tool is used** without (full) understanding
- Statistical tools applied because **supervisor/reviewer** asked
- Give up trying to understand, just use it as a **blackbox**



# Truth be told: everyone hates statistics

## Probability and statistics:

Difficult math, boring and totally useless, everyone hates it!



## Probability and statistics:

~~Difficult math, totally useless and so boring, everyone hates it!~~

- **Difficult:** mostly sums and products of fractions
- **Boring:** that's subjective, but yes, it may be boring
- **Useless:** definitely not! The basis of modern empiricism

# Truth be told: everyone hates statistics

## Probability and statistics:

~~Difficult math, totally useless and so boring, everyone hates it!~~

- Yes, we may hate it, but we also need it!
  - Knowing what we're doing can make us feel more at ease
  - It is worth the effort of overcoming initial resistance

## Probability and statistics:

- Yes, we may **hate** it, but we also **need** it!
  - Knowing what we're doing can make us feel more at ease
  - It is worth the effort of overcoming initial resistance

Ready? Let's go!

## Random variable: definition

- A **random variable** is a variable with no specific **value**
  - It takes some value within a (known) set of possible values
  - We are not interested in its actual value

Examples:

- A **human's age** takes values from 0 to 130 years
- **Water temperature** at sea level ranges from 0°C to 100°C
- A person's **handedness** can be right-handed, left-handed, both

Wooclap time!

## Random variable: examples

Are the following (interesting) random variables?

- 1. The **number of tentacles** of an octopus?

## Random variable: examples

Are the following (interesting) random variables?

- 1. The **number of tentacles** of an octopus?

→ No, always the same value

## Random variable: examples

Are the following (interesting) random variables?

- 2. The **distance** between the Earth and the Moon?

## Random variable: examples

Are the following (interesting) random variables?

- 2. The **distance** between the Earth and the Moon?  
→ Yes, it actually varies from 363K to 406K km

## Random variable: examples

Are the following (interesting) random variables?

- 3. A person's **vote** in the last presidential elections?

## Random variable: examples

Are the following (interesting) random variables?

- 3. A person's **vote** in the last presidential elections?  
→ Yes, the values are the candidates/parties running

## Random variable: examples

Are the following (interesting) random variables?

- 4. A person's **opinion** about how cute an octopus is?

## Random variable: examples

Are the following (interesting) random variables?

- 4. A person's **opinion** about how cute an octopus is?
  - No, ill-defined, no closed set of possible values
  - Actually, everyone finds them cute! ;-)

# Random variable: examples

Are the following (interesting) random variables?

- 1. The number of tentacles of an octopus? **No**
- 2. The distance between the Earth and the Moon? **Yes**
- 3. A person's vote in the last presidential elections? **Yes**
- 4. A person's opinion about how cute an octopus is? **No**

## In short

- A variable is not random if its value is fixed / constant
- We need to be able to describe its **set of possible values**
  - The set may be infinite (e.g. real numbers)

# Why do we need random variables?

---

- Use their **characteristics** to understand the data
- Model **features** and **evaluation metrics** as random variables
- Basic block in **probability** and **statistics**
  - People have been studying them for a while
  - Statistical tools associated to them can be useful

## Probability distributions

- Random variables come with probability distributions

$$P\{X = a\} = p(a) = 0.8 = 80\%$$

# Probability distributions

- Random variables come with probability distributions

$$P\{X = a\} = p(a) = 0.8 = 80\%$$

- $X$ : The random variable that we're interested in
- $a$ : The particular value of that random variable
- 0.8 (80%): The probability that variable  $X$  takes value  $a$

# Probability distributions

- Random variables come with probability distributions

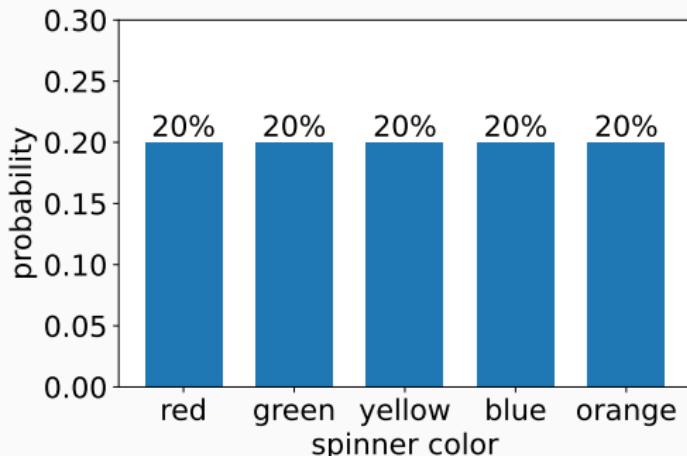
$$P\{X = a\} = p(a) = 0.8 = 80\%$$

- $X$ : The random variable that we're interested in
- $a$ : The particular value of that random variable
- 0.8 (80%): The probability that variable  $X$  takes value  $a$

Note: probabilities  $p(a)$  must sum up to 1 (for all values  $a$ )

## Example probability distributions

- $X_1$ : color of a 5-coloured spinner wheel

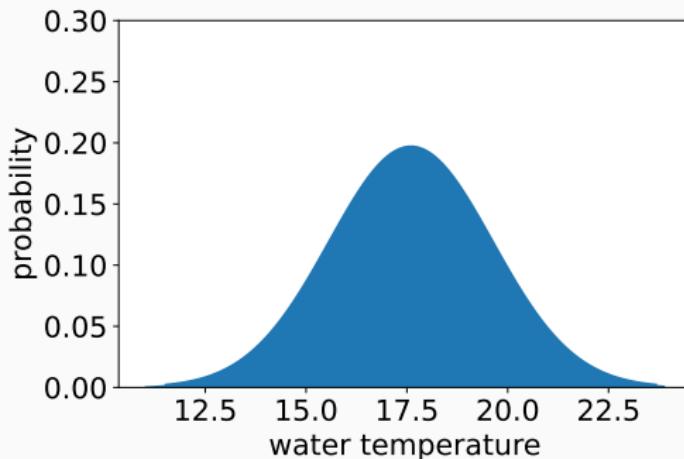


$$P\{X_1 = \text{red}\} = p(\text{green}) = \dots = p(\text{orange}) = \frac{1}{5}$$



## Example probability distributions

- $X_2$ : sea water temperature in July in Marseille



$$P\{X_3 < 17.6\} = 0.5$$

Wooclap time!

# Probability distribution or not?

Which of the following are proper probability distributions? Why?

a)

$x_i$	$p(x_i)$
1	0.4
2	-0.2
3	0.8

b)

$x_i$	$p(x_i)$
0.4	0.4
0.35	0.35
0.25	0.25

c)

$x_i$	$p(x_i)$
-1	0.4
-2	0.2
-3	0.8

d)

$x_i$	$p(x_i)$
-1	0.4
0	0.2
1	0.2
2	0.1

# Probability distribution or not?

Which of the following are proper probability distributions? Why?

a)

$x_i$	$p(x_i)$
1	0.4
2	-0.2
3	0.8

No,  $p(2) < 0$

b)

$x_i$	$p(x_i)$
0.4	0.4
0.35	0.35
0.25	0.25

Yes, sum=1

c)

$x_i$	$p(x_i)$
-1	0.4
-2	0.2
-3	0.8

No, sum > 1

d)

$x_i$	$p(x_i)$
-1	0.4
0	0.2
1	0.2
2	0.1

No, sum < 1

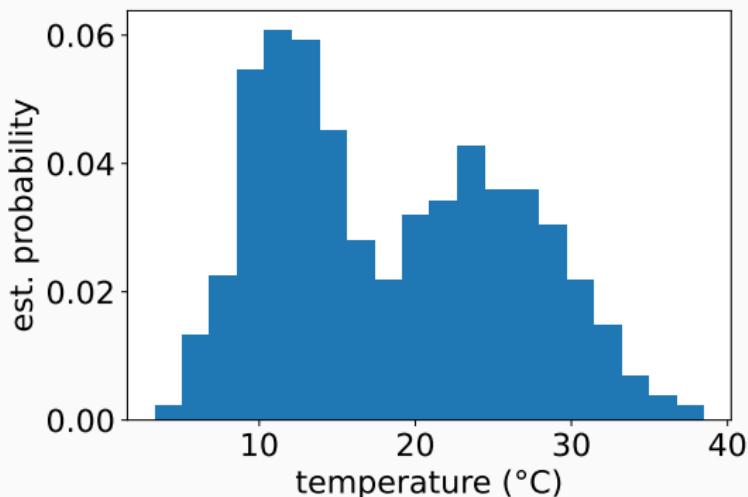
# From probabilities to statistics

- Probability distributions are theoretical **abstractions**
  - We often learn probabilities with toy examples
  - In practice,  $X$ 's "real" distribution is not accessible
- A **sample** is often used to **estimate** the probabilities
  - Most of the time, probabilities are approximated
  - **Proportion in sample (%)** → estimated probability

$$\frac{\text{count}(a)}{n} \approx P\{X = a\}$$

## Random samples

- Randomly collect a finite set of data points
- Example: temperature of a sensor in a power plant  
→ Size: 365 days → [10.1, 14.0, 8.9, ..., 12.5, 15.3, 13.3]



Estimated probability distribution = normalized histogram

# Sampling: example

## Jupyter notebook 1 & 2

1. Open the dataset using `pandas.read_csv()`
2. Explore the different columns and their values
3. Make a histogram of the compositionality column  
→ This is an estimate of its distribution!

## Compositionality dataset

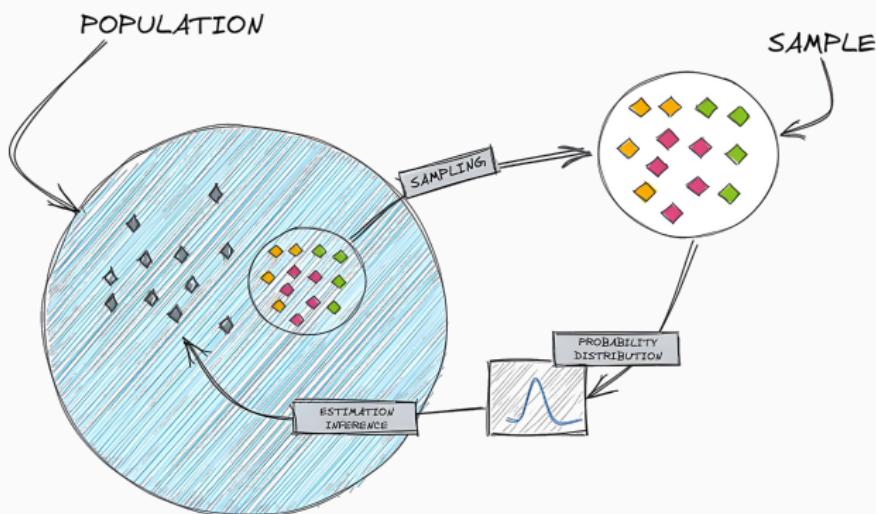
- Is a dry run literally a run which is dry?  
→ not at all ← 0 - 1 - 2 - 3 - 4 - 5 → absolutely yes
- Compositionality score: average rating of 10-15 annotators
- Sample: 180 compounds in French

Source: <https://aclanthology.org/J19-1001/>



# Why do we need samples?

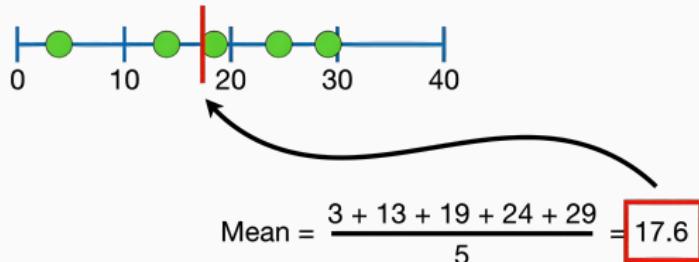
- A **representative sample** can inform us about the whole
  - Full data not available: **infer** properties of (unknown) distribution
  - Draw **generalisable conclusions** in the presence of uncertainty



# Sample mean / average

- A single value at the **center** of the sample  
→ Summarise the whole data set with a single number

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{1}{n} \sum_{i=1}^n x_i$$



Wooclap time!

## Mean / average quiz

---

- Is the mean a probability (value between 0 and 1)?

## Mean / average quiz

- Is the mean a probability (value between 0 and 1)?  
→ No, it depends on the values (arbitrary range)
- Is the value of the mean contained in the sample?

## Mean / average quiz

- Is the mean a probability (value between 0 and 1)?  
→ No, it depends on the values (arbitrary range)
- Is the value of the mean contained in the sample?  
→ No, it can be a new value, not contained in the sample
- Is the value of the mean always positive?

# Mean / average quiz

- Is the mean a probability (value between 0 and 1)?  
→ No, it depends on the values (arbitrary range)
- Is the value of the mean contained in the sample?  
→ No, it can be a new value, not contained in the sample
- Is the value of the mean always positive?  
→ No, e.g. if the variable only takes negative values

# Data dispersion

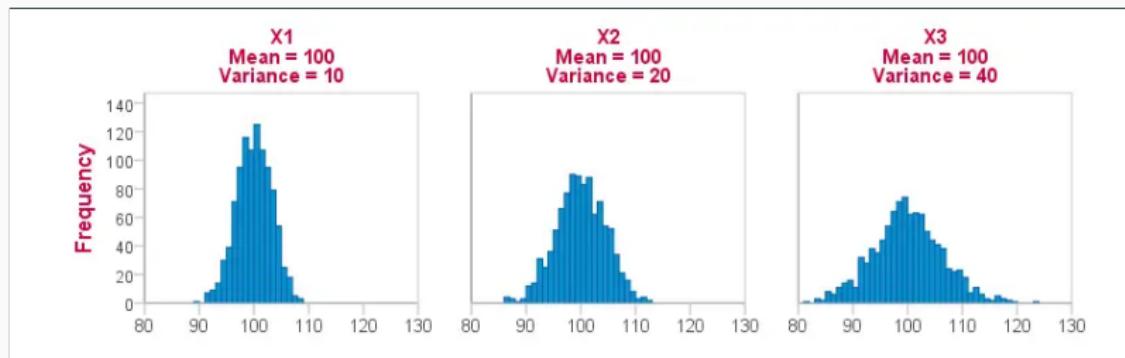
- Mean does not take into account **data dispersion**

$$S_1 \quad [0, 0, 0, 0] \quad \overline{S_1} = 0$$

$$S_2 \quad [-4, -4, 4, 4] \quad \overline{S_2} = 0$$

$$S_3 \quad [-6, -2, 1, 7] \quad \overline{S_3} = 0$$

$$S_4 \quad [-15000, 15000] \quad \overline{S_4} = 0$$



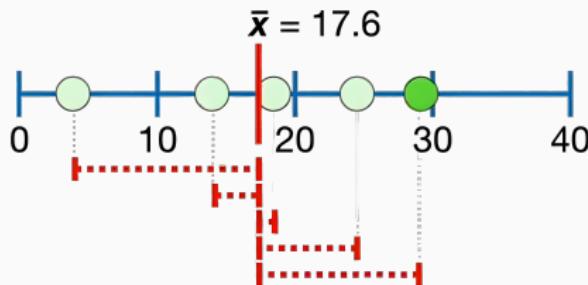
Source: <https://www.spss-tutorials.com/descriptive-statistics-one-metric-variable/>

# Variance

- Variance characterises the dispersion/spread of a distribution
  - Intuition: average distance from the mean
  - $(x_i - \bar{x})$  can be positive or negative  $\Rightarrow$  square it!

$$\text{Var}(X) = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1}$$

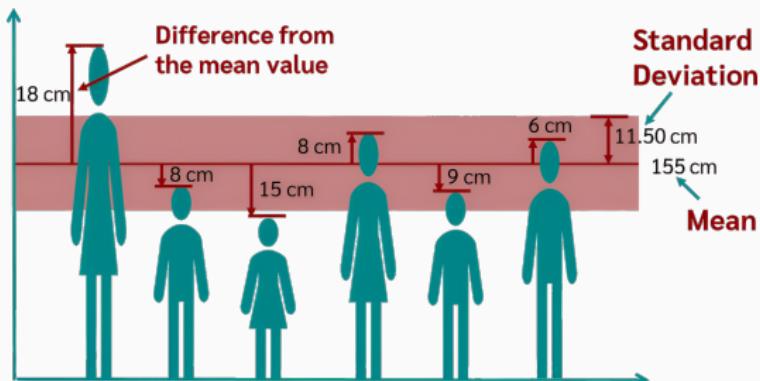
→ Variance is always positive, differently from mean



# Standard deviation

- Variance averages *squared* differences
  - Its absolute value is hard to interpret
  - Bring back to original value range → squared root
- The squared root of variance is called **standard deviation**

$$\sigma = \sqrt{\text{Var}(X)}$$



# Estimated standard deviation

- Population standard deviation:

$$\sigma_x = \sqrt{\sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n}}$$

- Sample standard deviation, unbiased estimator:

$$s_x = \sqrt{\sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1}}$$

- Why? <https://www.youtube.com/watch?v=sHRBg6BhKjI>

In practice, we only need  $s_x$  → Most stats libraries' default

# Calculating mean and standard deviation

## Jupyter notebook 3

1. Open dataset containing 180 compositionality scores
2. Use Pandas' `comp.describe()` to obtain a summary
3. Is the obtained standard deviation  $\sigma_X$  or  $s_X$ ?

# One distribution to rule them all

The **Normal** distribution

$$P\{a < X < b\} = \int_a^b \frac{1}{\sqrt{2\pi\sigma^2}} \exp^{\frac{x-\mu}{\sigma}}$$

# One distribution to rule them all

The **Normal** distribution

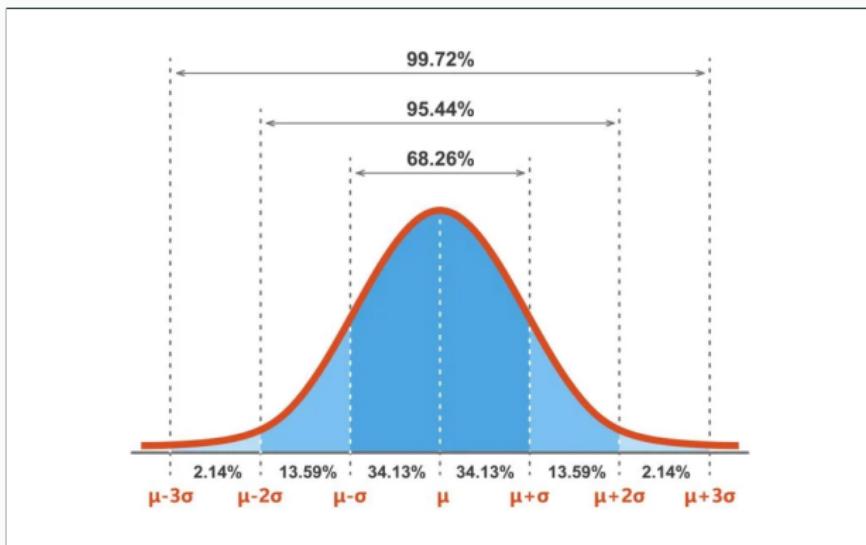
$$P\{a < X < b\} = \int_a^b \frac{1}{\sqrt{2\pi\sigma^2}} \exp^{\frac{x-\mu}{\sigma}}$$

Who cares!

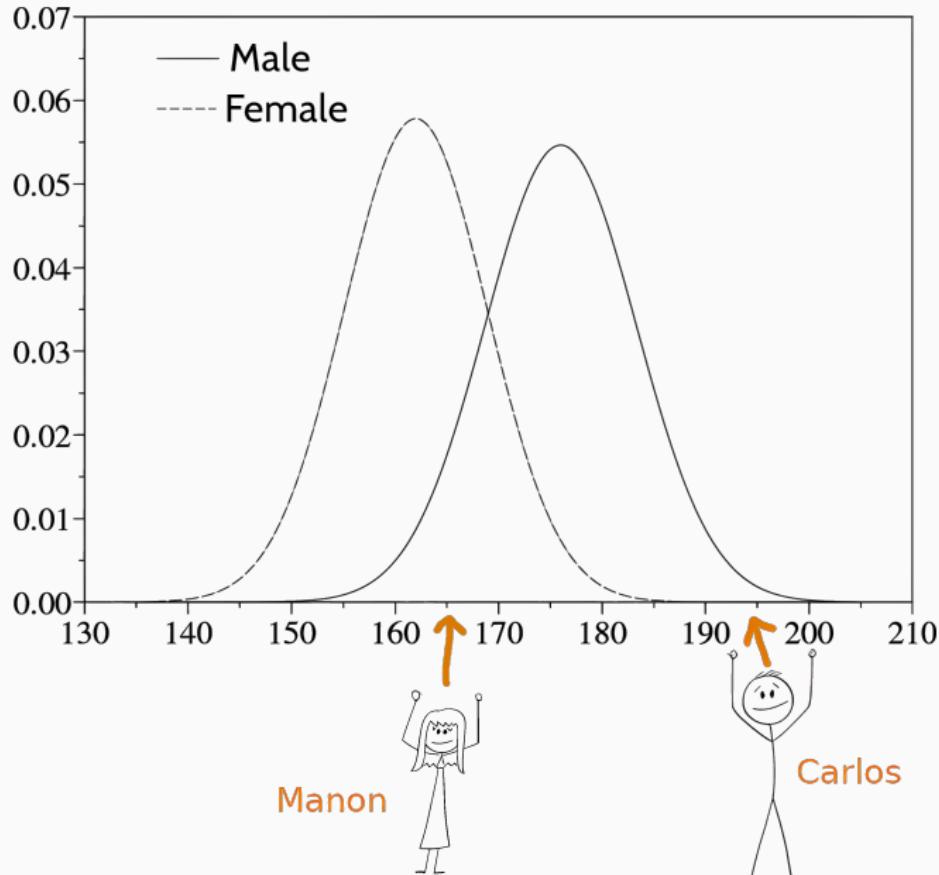
# One distribution to rule them all

## The **Normal** distribution

- Probability density function is a symmetric **bell-shaped curve**
- Specified by mean  $\mu$  (center) and std. deviation  $\sigma$  (wideness)  
→ 99.7% of probability between  $\mu - 3\sigma$  and  $\mu + 3\sigma$



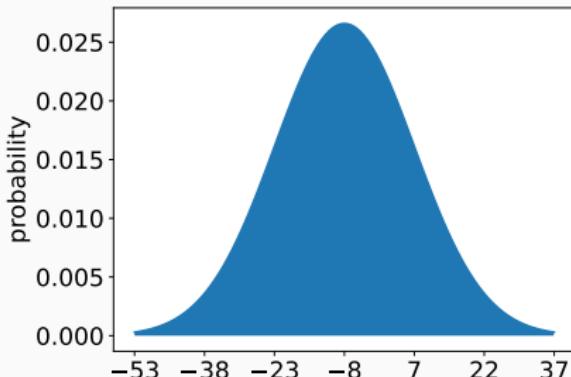
## Normal distribution: example



Wooclap time!

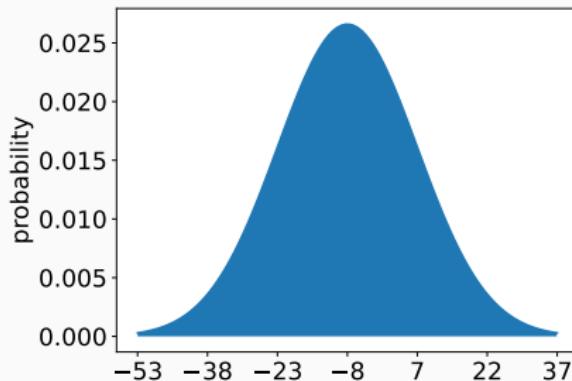
# Who's that normal?

1. What are the  $\mu$  and  $\sigma$  parameters for the following curve?



# Who's that normal?

- What are the  $\mu$  and  $\sigma$  parameters for the following curve?

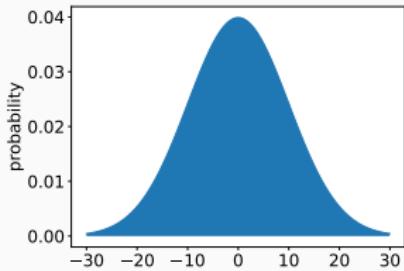


$$\mu = -8 \text{ and } \sigma = 15$$

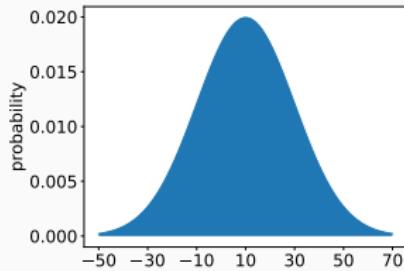
# Who's that normal?

1. What are the  $\mu$  and  $\sigma$  parameters for the following curve?
2. Which curve corresponds to  $\mu = 10$  and  $\sigma = 20$ ?

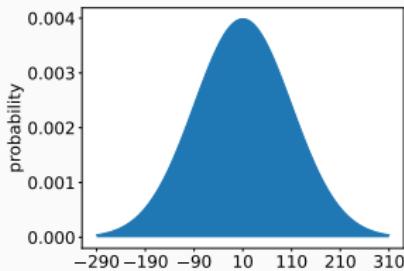
a)



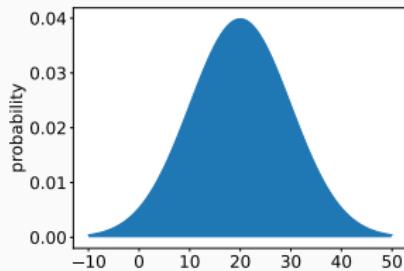
b)



c)

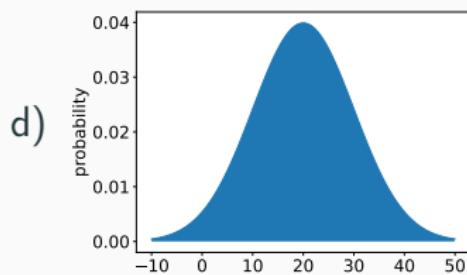
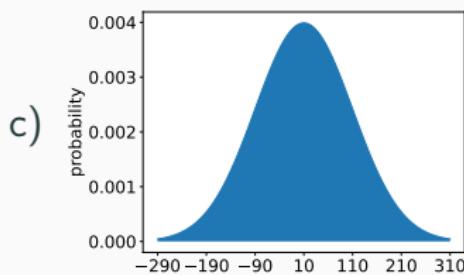
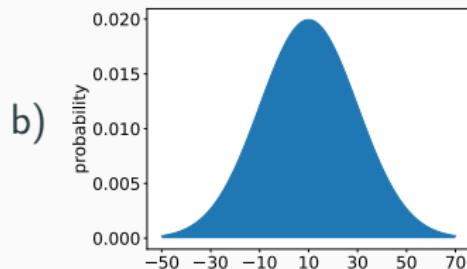
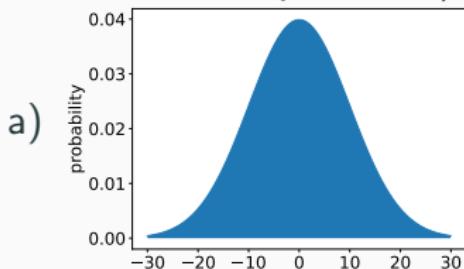


d)



# Who's that normal?

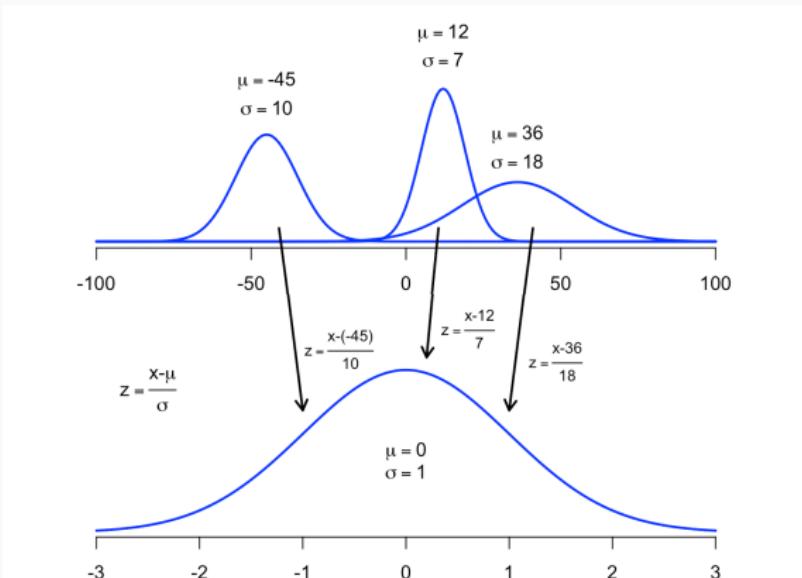
1. What are the  $\mu$  and  $\sigma$  parameters for the following curve?
2. Which curve corresponds to  $\mu = 10$  and  $\sigma = 20$ ?



curve b) – notice different heights

# Standardization

- Calculate probability → integration ( $\langle \circ \rangle$  aaaaah!)  
→ Normal is impossible to integrate analytically
- In practice:  
→ Standardize  $z = \frac{x-\mu}{\sigma}$ , then **lookup table** of  $\Phi(a)$



Wooclap time!

# The most famous probability distribution

Why is the normal distribution so important?

# The most famous probability distribution

Why is the normal distribution so important?

- Turns out **most measurements** are normally distributed
- Used in many statistical tools, e.g. **hypothesis testing**
- Plays a central role in describing **estimated means**

# It's normal to be average

- Normalised sum of random variables is **normally distributed**<sup>2</sup>  
→ Even if the variables are **not** normally distributed!
- The **mean  $\bar{X}$**  of a sample is normally distributed  
→ Comes in handy to analyse **averages**
- This is known as the **central limit theorem** (CLT)

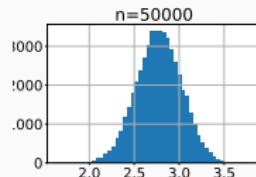
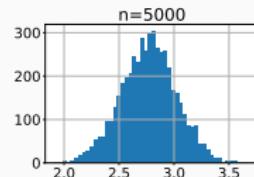
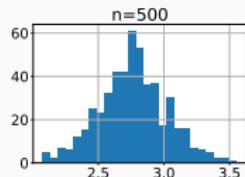
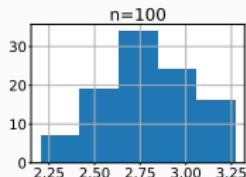
---

<sup>2</sup>Actually, independent and identically distributed (i.i.d.) random variables.

# Central limit theorem: example

Jupyter notebook 4 & 5

1. Build  $n$  random samples of size 30 from compositionality data
2. Calculate mean of each random sample, save values
3. Estimate sample mean's distribution with histogram  
→ What happens when  $n$  increases?



## In short

- Random variables and probability distributions
  - Theoretical model for features and metrics
  - In practice, estimated using sampling
- Mean and standard deviation characterise the data
- Normal distribution: bell shaped around the mean
  - Useful to characterise values that are means (CLT)

## In short

- Random variables and probability distributions
  - Theoretical model for features and metrics
  - In practice, estimated using sampling
- Mean and standard deviation characterise the data
- Normal distribution: bell shaped around the mean
  - Useful to characterise values that are means (CLT)

Now we're ready for the next steps!



# Outline

---

Statistics in a nutshell

Evaluation metrics

Statistical significance

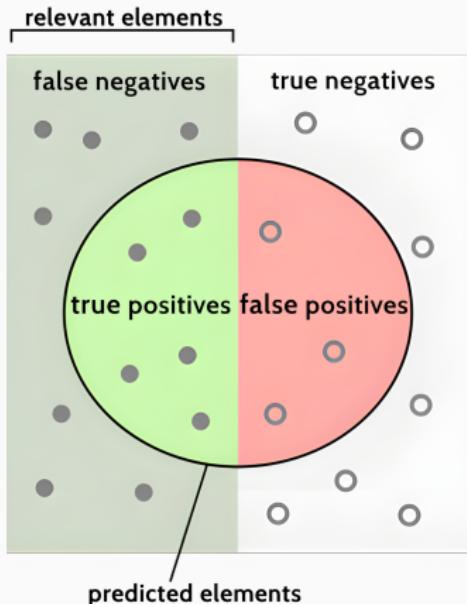
Discussion

## Disclaimer: all metrics are incomplete

- Ideally: measure a hidden variable or phenomenon
- In practice: measure what we can observe
  - Formulation is simple enough to be interpretable
- Metrics are partial views of the results

# Classification framework

- ***tp*: true positives**  
→ Correctly predict as positive
- ***tn*: true negatives**  
→ Correctly predict as negative
- ***fp*: false positives**  
→ Predict positive, should be negative
- ***fn*: false negatives**  
→ Predict negative, should be positive



Source: image adapted from Wikipedia

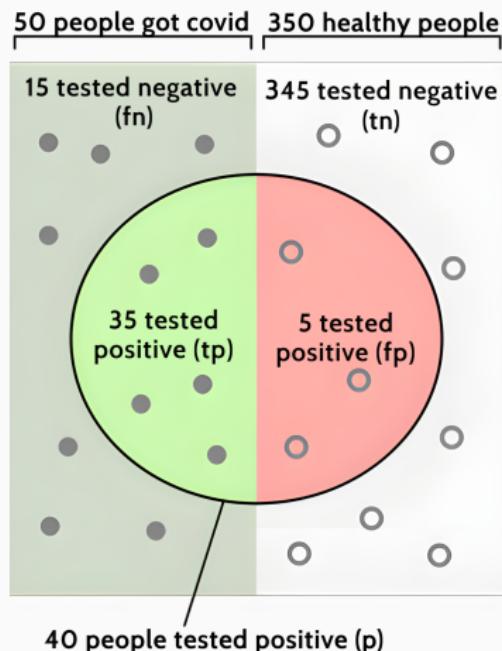
## Classification framework: example

A group of 400 people did a covid-19 test. 50 people really got covid (relevant), the other 350 people do not. The test is positive for 40 people, out of which 35 really have covid.

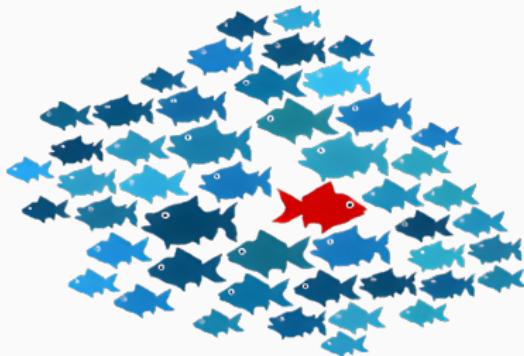
## Classification framework: example

A group of 400 people did a covid-19 test. 50 people really got covid (relevant), the other 350 people do not. The test is positive for 40 people, out of which 35 really have covid.

- $p: 40$
- $n: 360$
- $tp: 35$
- $tn: 345$
- $fp: 5$
- $fn: 15$



# Accuracy



$$Accuracy = \frac{tp+tn}{tp+tn+fp+fn}$$

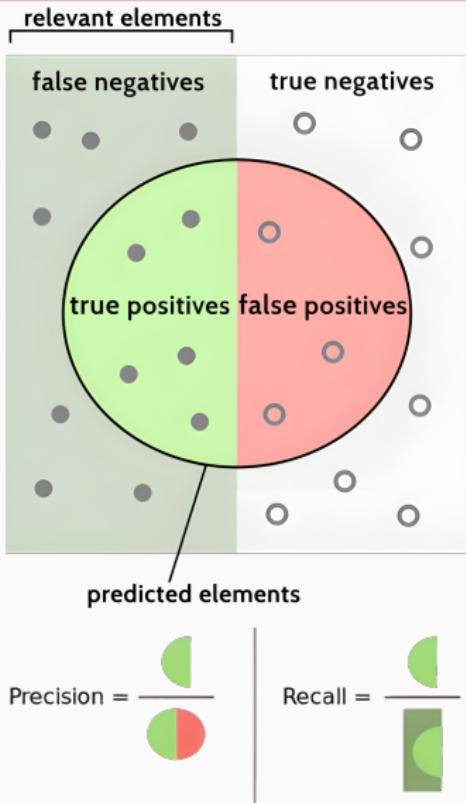
- Percentage of **well classified** items
- Incomplete description of the method's performance

Source: Image: Devin Soni, [towardsdatascience.com](https://towardsdatascience.com)

# Precision, recall, F-score

- True positive ratios:
  - Precision  
 $\rightarrow tp/(tp + fp)$
  - Recall = Sensitivity  
 $\rightarrow tp/(tp + fn)$
  - Specificity:  
 $\rightarrow tn/(tn + fp)$
- Complementary measures, report both
- F-score: prec. & recall harmonic mean

$$F = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$



Source: image adapted from Wikipedia

# Multi-class classification

- Precision, recall and F-score are calculated **for each class!**
- For the case of 2 classes, often 1 is more important
  - Spam detection: we're interested in the SPAM class
- Multi-class classification with  $k$  classes:
  - $k$  precision scores
  - $k$  recall scores
  - $k$  F-scores

Wooclap time!

## Métriques de classification: exercice

Suppose you are working on a spam detection system for email. You have developed a machine learning model that has been evaluated on a test set containing 1,000 emails, 200 of which were really spam. The model classified 150 of these 200 emails as spam. 25 emails were classified as spam when they should not have been.

- Obtain the number of tp, tn, fp and fn.
- Calculate precision, recall, and accuracy.

## Métriques de classification: exercice

Suppose you are working on a spam detection system for email. You have developed a machine learning model that has been evaluated on a test set containing 1,000 emails, 200 of which were really spam. The model classified 150 of these 200 emails as spam. 25 emails were classified as spam when they should not have been.

- Obtain the number of tp, tn, fp and fn.
  - Calculate precision, recall, and accuracy.
- 
- tp: 150
  - tn: 775
  - fp: 25
  - fn: 50
- Precision:  $\frac{150}{150+25} = 0.857$
  - Recall:  $\frac{150}{150+50} = 0.75$
  - Accuracy:  $\frac{150+775}{1000} = 0.925$

# Accuracy and class imbalance

- Example: hate speech detection in tweets
  - Only a small percentage ( $\sim 1\%$ ) are hateful
  - Let us annotate **everything as not hateful**
  - My model has an **accuracy** of 99%! So powerful!



# Accuracy and class imbalance

- Example: hate speech detection in tweets
  - Only a small percentage ( $\sim 1\%$ ) are hateful
  - Let us annotate everything as not hateful
  - My model has an accuracy of 99%! So powerful!



High class imbalance: accuracy is misleading/not relevant

## F-score or F-measure

- F-score (or F-measure): harmonic mean of precision and recall

$$F = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

- F-score can be weighted to favour precision or recall
  - $\beta=0.5$ : More weight on precision, less weight on recall
  - $\beta=1$ : Balance the weight on precision and recall
  - $\beta=2$ : Less weight on precision, more weight on recall

$$F_\beta = (1 + \beta^2) \times \frac{\text{Precision} \times \text{Recall}}{\beta^2 \text{Precision} + \text{Recall}}$$

## Other metrics (see backup slides)

- ROC curve / Area under the curve
  - Real prediction, threshold
- (Mean) average precision
  - Real prediction, binary gold classes
- Structured prediction
  - Compare trees, graphs, clusters...
- ...

## Warning!

⚠️ BIAS ALERT ⚠️

## Metric bias

- Choosing the most interesting metrics **after** the experiment
  - E.g. a few positive examples, thousands of negative examples
  - Report accuracy, omit f-score: strongly affected by class imbalance
- Evaluation metrics must be defined **before** the experiment

*“Just because statistical software [...] generates output by default does not mean that all output is useful or relevant.”*

Source: Bruce et al (2020)

# Goodhart's law

"When a measure becomes a target, it ceases to be a good measure"

- Cobra effect
- Reinforcement learning policy
- Grade-oriented education system
- Risk: optimise evaluation metric **at any expense**
  - Overfitting, low generalisation
  - Forgetting the research question
  - Frustration with unrealistic goals
  - ...

Source: Thanks to François Hamonic

# Evaluation metrics and statistics

- Evaluation **metrics** are calculated on **datasets**
  - **Datasets** = samples
  - **Metrics** = random variables

- Evaluation **metrics** are calculated on **datasets**
  - **Datasets** = **samples**
  - **Metrics** = **random variables**

## Questions

1. How can we compare these **random variables**?
2. Can we trust that **sample**-based estimations generalise?

# Outline

---

Statistics in a nutshell

Evaluation metrics

Statistical significance

Discussion

# Year 3000...

The Earth is finally a **safe and pleasant** place for humans again.

However, 1000 years of global warming unleash a **dangerous permafrost bacteria**.

The bacteria starts to **infect human** hosts, causing a mysterious disease.

Centuries in insipid watery ice made the bacteria **obsessive** about...



# ...vanilla ice-cream! ❤



The illness is called

- Compulsive
- Obsessive
- Vanilla
- Ice-cream
- Disease



# CHAOS !!

The bacteria spreads rapidly, and infected humans start **eating tons of vanilla ice-cream**.

Milk prices rise to the stratosphere, ice-cream makers strike, diabetes and obesity break records...

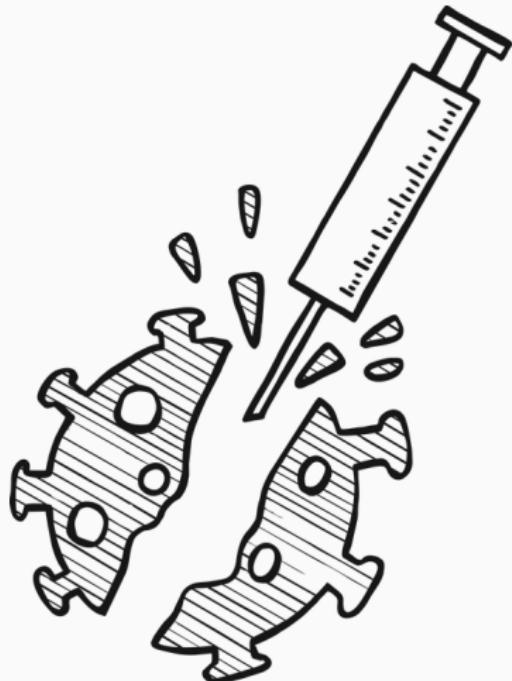
Governments impose ice-cream lockdowns, interplanetary travel is forbidden, **panic** everywhere!



After months of an unprecedented crisis...

A lab finally announces a **vaccine**  
at phase 3 of evaluation!

This requires an experiment called  
**randomized control trial**



# Randomized control trial

Group A  
Vaccine

Group B  
Placebo

Conclusion:

The vaccine works!

What a relief for humanity!

Average nb. ice-creams/day (ICD):

- Group A:  $ICD_A = 1.47$
- Group B:  $ICD_B = 1.56$

$$ICD_A(\text{vaccine}) < ICD_B(\text{placebo})$$



## But... maybe humans forgot all about statistics?

- Is the observed difference large enough?
  - $ICD_A = 1.47$  ice/creams per day
  - $ICD_B = 1.56$  ice/creams per day

$$\delta = ICD_B - ICD_A = 0.09$$

- Maybe this result is due to **randomness in sampling**
  - Affects our conclusion that vaccine (A) better than placebo (B)?

## But... maybe humans forgot all about statistics?

- Is the observed difference large enough?
  - $ICD_A = 1.47$  ice/creams per day
  - $ICD_B = 1.56$  ice/creams per day

$$\delta = ICD_B - ICD_A = 0.09$$

- Maybe this result is due to **randomness in sampling**
  - Affects our conclusion that vaccine (A) better than placebo (B)?

Given the samples, the metrics, and the experiment's conditions:

What is the probability of making a **false claim** if we conclude that *A* better than *B* in general?

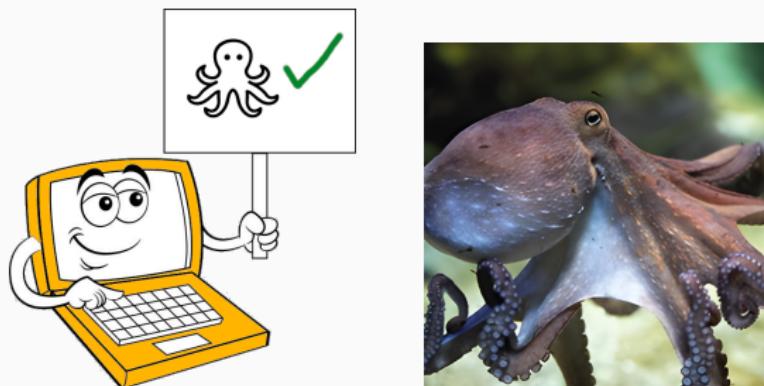
## System comparison: example

- Our Baseline system classifies images  
→ Two categories: octopus or not octopus



## System comparison: example

- Our Baseline system classifies images  
→ Two categories: octopus or not octopus



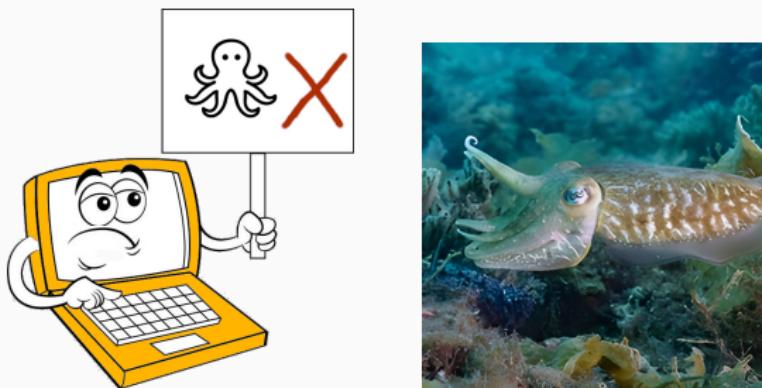
## System comparison: example

- Our Baseline system classifies images  
→ Two categories: octopus or not octopus



## System comparison: example

- Our Baseline system classifies images  
→ Two categories: octopus or not octopus



## System comparison: example

- Our Baseline system classifies images  
→ Two categories: octopus or not octopus



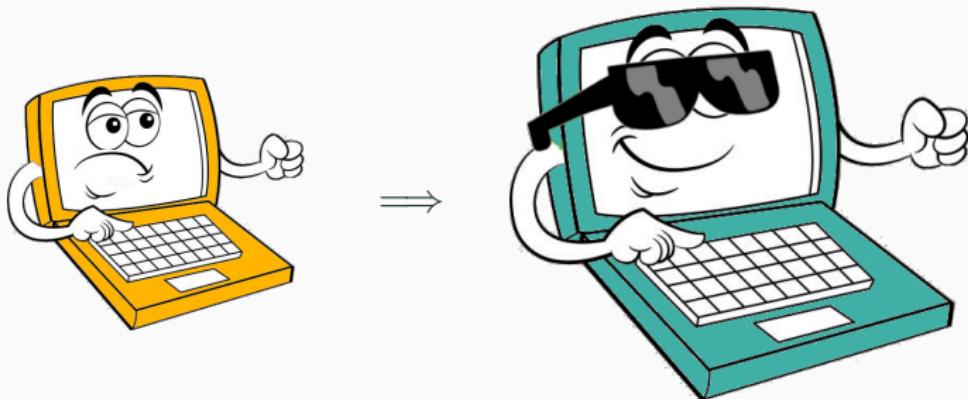
# System comparison: example

- Our **Baseline** system classifies images
  - Two categories: octopus or not octopus
- Sometimes it makes **mistakes**



## System comparison: example

- We developed an Awesome new system!
  - E.g. the new system was trained on more data



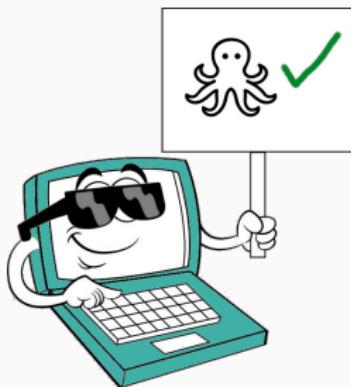
# System comparison: example

- We developed an Awesome new system!  
→ E.g. the new system was trained on more data



# System comparison: example

- We developed an **A**wesome new system!  
→ E.g. the new system was trained on more data
- It seems that it makes less **mistakes** ⇒ 🎉



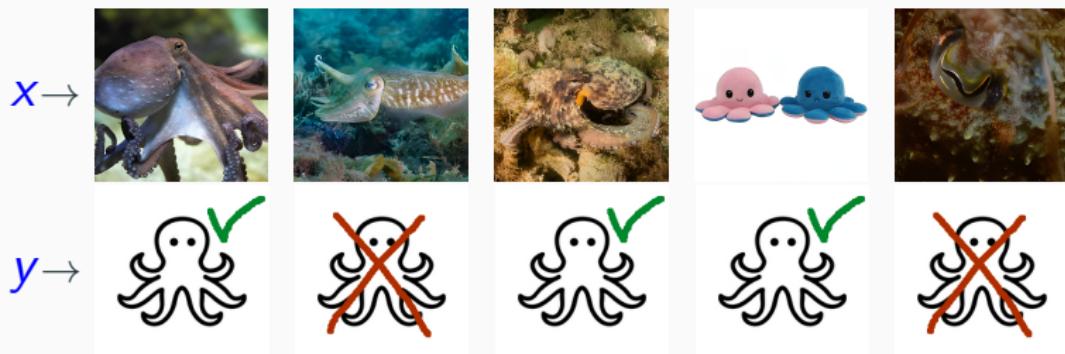
## Test set

---

- Is A really better than B?
  - Testing on a few examples is not enough!
- Use a test set containing (x,y) pairs
  - x - sea animal images
  - y - reference octopus/other labels
- The test set was not used to develop the system

## Test set: example

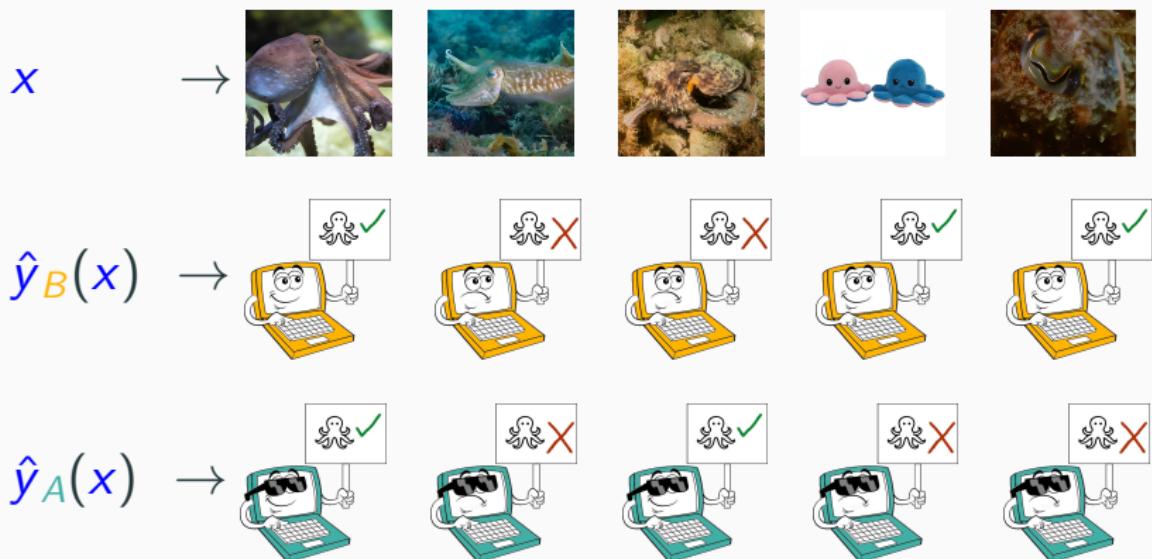
Images  $x$  selected to be in the held-out test set



Gold labels  $y$  considered as reference (e.g. annotated by humans)

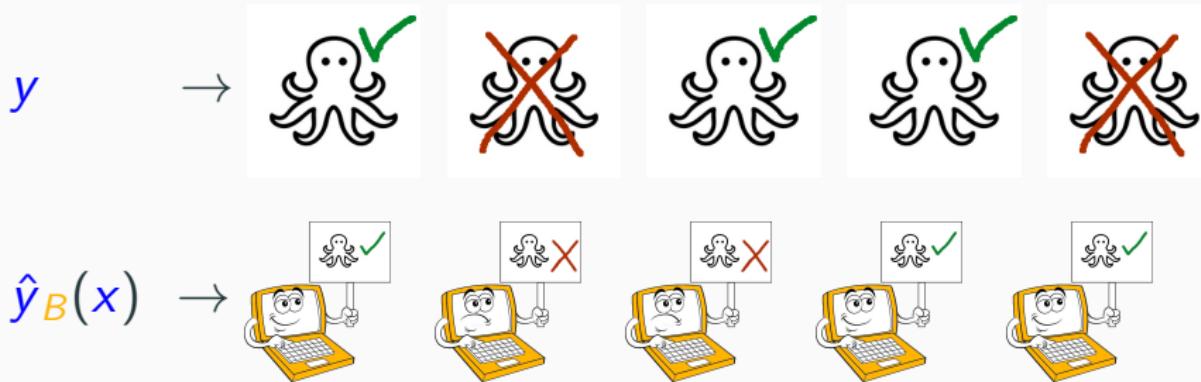
# System predictions

Both systems generate predictions  $\hat{y}$  for test set instances  $x$



## ⚡ Evaluation metrics ⚡

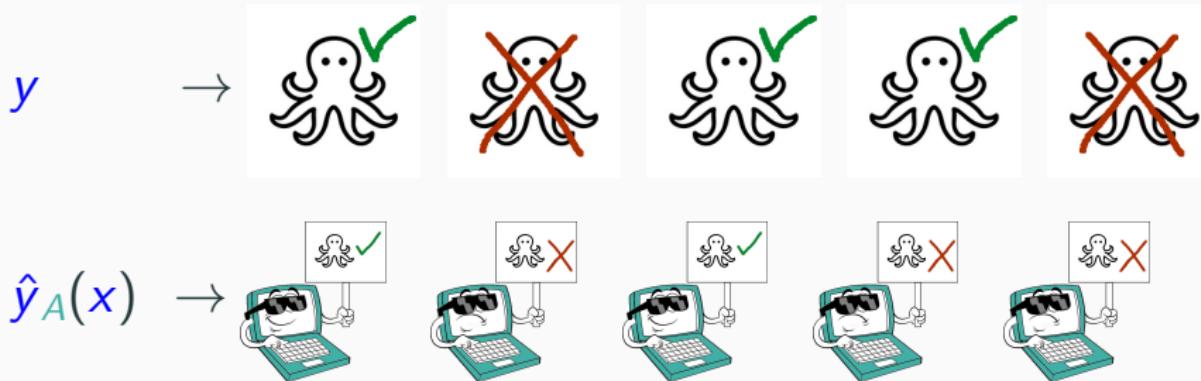
Compare predictions  $\hat{y}_B$  and  $\hat{y}_A$  to reference  $y$



$$M(B, x, y) = \frac{3}{5} = 0.6$$

## ⚡ Evaluation metrics ⚡

Compare predictions  $\hat{y}_B$  and  $\hat{y}_A$  to reference  $y$



$$M(A, x, y) = \frac{4}{5} = 0.8$$

Wooclap time!

## System score comparison

- The accuracies of both systems are:

$$M(B, x, y) = \frac{3}{5} = 0.6$$

$$M(A, x, y) = \frac{4}{5} = 0.8$$

- It seems like A is better than B
- The difference (delta) is positive

$$\delta_{A-B}(x, y) = M(A, x, y) - M(B, x, y) = 0.8 - 0.6 = 0.2$$

## System comparison: example

We obtained a much larger test set  $x',y'$



We compare A and B again and obtain:

## System comparison: example

We obtained a much larger test set  $x', y'$



We compare A and B again and obtain:

$$\begin{aligned}\delta_{A-B}(x', y') &= M(A, x', y') - M(B, x', y') \\ &= 0.7612 - 0.7586 \\ &= \boxed{0.0026}\end{aligned}$$

## System comparison: example

We obtained a much larger test set  $x', y'$



We compare **A** and **B** again and obtain:

$$\begin{aligned}\delta_{A-B}(x', y') &= M(A, x', y') - M(B, x', y') \\ &= 0.7612 - 0.7586 \\ &= \boxed{0.0026}\end{aligned}$$

- Can we still affirm that **A** is better than **B**?
- If we add or remove a couple of images, could the result flip?

# Interpreting delta

$$\delta_{A-B}(x, y) = M(A, x, y) - M(B, x, y)$$

- Delta: compare both systems with a single value
  - A better than B →  $\delta_{A-B}(x, y) > 0$
  - A equivalent to B →  $\delta_{A-B}(x, y) = 0$
  - A worse<sup>3</sup> than B →  $\delta_{A-B}(x, y) < 0$
- In some disciplines,  $\delta_{A-B}(x, y)$  is called effect

---

<sup>3</sup>Yes, the old Baseline may outperform the new Awesome system!

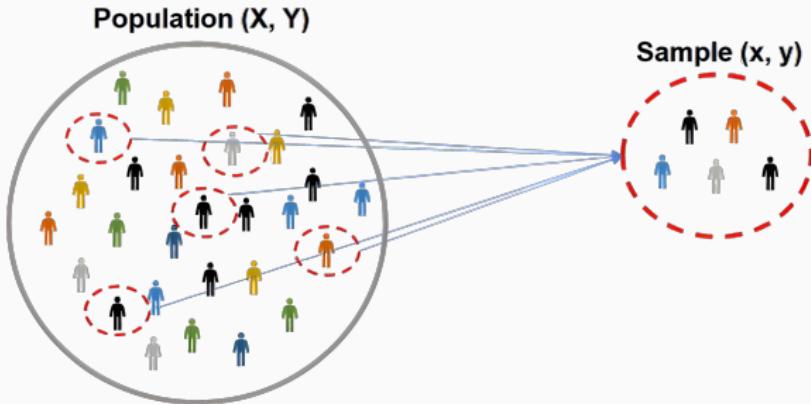
# System comparison in a nutshell

1. We develop a system  $A$  supposed to be better than  $B$
2. To verify this, we apply both systems to the same test set:
  - Get output of system  $A$  on the test set  $(x, y)$
  - Get output of system  $B$  on the test set  $(x, y)$
3. Calculate the evaluation metric  $M(\cdot)$  for both outputs

$$\delta_{A-B}(x, y) = M(A, x, y) - M(B, x, y)$$

4. Large positive  $\delta_{A-B}(x, y) \Rightarrow$
5. In practice,  $\delta_{A-B}(x, y)$  is often small 😊

# Test sets as random samples

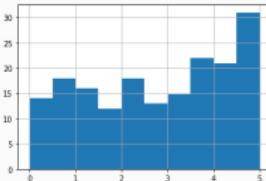


- Could the observed  $\delta_{A-B}(x, y) > 0$  be due to **sampling error** ?
  - $(x, y)$  is a sample of joint random variables  $(X, Y)$
  - What difference/effect would we get for another sample  $(x', y')$  ?
- What is the probability that  $A$  is actually no better than  $B$ 
  - If we ever had access to the “real” distribution of  $(X, Y)$ ?

## Effects as random variables

- We obtain a single  $\delta_{A-B}(x, y)$  value
- This value depends on the test set  $(x, y)$ , which is a sample
- We see  $\delta_{A-B}(x, y)$  as a sampled value of a **random variable**

$$\delta_{A-B}(X, Y) \sim$$



## P-value

---

- **P-value:** probability of obtaining at least  $\delta_{A-B}(x, y)$ 
  - When in reality, A is no better than B

- **P-value:** probability of obtaining at least  $\delta_{A-B}(x, y)$ 
  - When in reality, A is no better than B

## Warning!

- p-value  $\neq$  probability that your conclusion is wrong!
- Rather, probability that chance is the only reason for the delta

Wooclap time!

# P-value: example



We have one delta value obtained on the large dataset ( $x', y'$ )

$$\delta_{A-B}(x', y') = 0.0026$$

# P-value: example



We have one delta value obtained on the large dataset ( $x', y'$ )

$$\delta_{A-B}(x', y') = 0.0026$$

If we had all possible images of sea creatures and their classes

→ Imagine we have access to the real distribution  $\delta_{A-B}(X, Y)$

- P-value = probability of obtaining 0.0026 difference (or more)
- If A is actually no better than B

# Hypothesis testing

- $H_0: \delta_{A-B}(X, Y) \leq 0 \implies$  if true, then  $A$  not better than  $B$
- $H_1: \delta_{A-B}(X, Y) > 0$
- Goal: reject  $H_0$ 
  - Conclusion: **significant** difference between the systems

## Hypotheses

- $H_0: \delta_{A-B}(X, Y) \leq 0$
  - $H_1: \delta_{A-B}(X, Y) > 0$
- 
- **P-value:** probability of observing  $\delta_{A-B}(x, y)$  while  $H_0$  true  
→ Intuition: if  $H_0$  is true, large  $\delta_{A-B}(x, y)$  is unlikely

# Hypothesis testing: example



$$\text{p-value} = P\{\delta_{A-B}(X, Y) \geq 0.0026\} \text{ assuming } \delta_{A-B}(X, Y) \leq 0$$

Usually, if p-value small enough  $\implies$  we conclude A better than B

# Type I errors

- Type I error: **false positive**  
→ Rejecting  $H_0$  when it is actually true

**Conclusion** of the test:



is better than

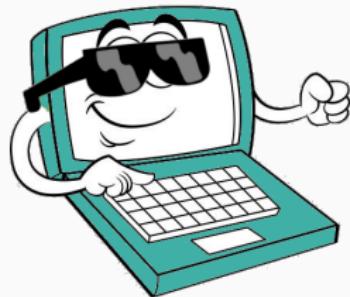


**Reality:** But it isn't better !

## Type II errors

- Type II error: **false negative**  
→ Not rejecting  $H_0$  when it is actually false

**Conclusion** of the test:



is not better than



**Reality:** But it is better !

# Goal

- Probability of type-I error is upper bounded by  $\alpha$   
→  $\alpha$  is called the **significance level** or threshold
- Probability of type-II error is as low as possible  
→ Test **power**: ability to avoid type-II errors



imgflip.com

## Statistically significant result

$p\text{-value} < \alpha \implies \text{statistically significant!}$  

- p-value: probability of extreme outcome under null hypothesis
- $\alpha$ : significance level (threshold)
  - Usual “magic” value:  $\alpha = 0.05$

# Statistically significant result

$p\text{-value} < \alpha \implies \text{statistically significant!}$  

- p-value: probability of extreme outcome under null hypothesis
- $\alpha$ : significance level (threshold)
  - Usual “magic” value:  $\alpha = 0.05$

The word **significant** should not be used for anything else

Wooclap time!

# What's a “good” p-value?

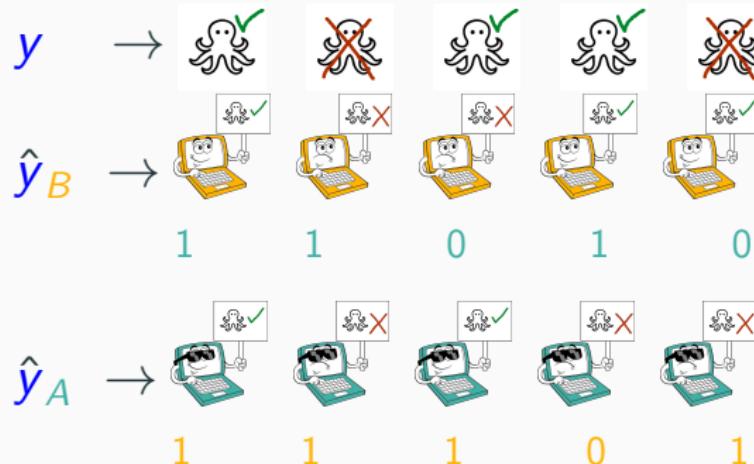
---

- We want the p-value to be **below**  $\alpha$  (significance level)
  - Usual threshold:  $\alpha = 0.05$
  - **p-value=0.045** and **p-value=0** are below the  $\alpha$  threshold
- If so, we say that the **A-B** difference is **significant**

# How can we estimate p-values?

- P-value depends on  $\delta_{A-B}(X, Y)$  probability distribution
- Which in turn depends on  $M(A, x, y)$  and  $M(B, x, y)$ 
  - Remember:  $M(\cdot)$  is our evaluation metric
- $M(\cdot)$ 's distribution determines that of  $\delta$ 
  - Study the probability distribution of evaluation metric  $M(\cdot)$

# Accuracy is an average

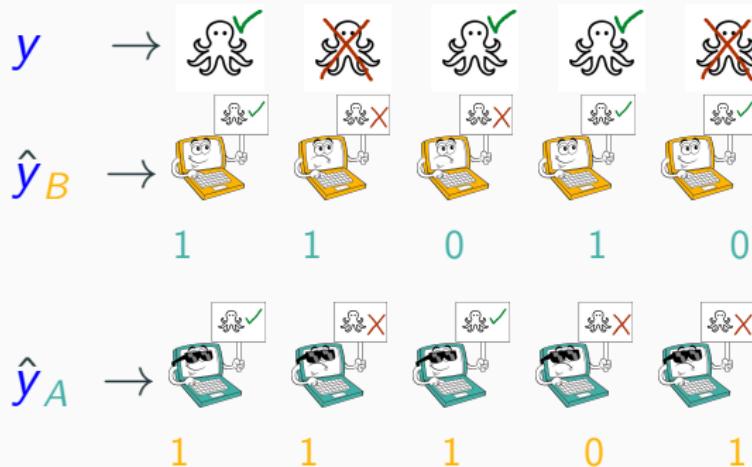


$$Acc_B = \frac{1+1+0+1+0}{5} = \frac{3}{5}$$

$$Acc_A = \frac{1+1+1+0+1}{5} = \frac{4}{5}$$

Accuracy is an average

# Accuracy is an average



$$Acc_B = \frac{1+1+0+1+0}{5} = \frac{3}{5}$$

$$Acc_A = \frac{1+1+1+0+1}{5} = \frac{4}{5}$$

Accuracy is an average  $\implies$  normally distributed!

$\rightarrow$  Thanks to the central limit theorem

# The t-test for paired samples

- T-test: hypothesis testing for **normally distributed variables**  
→ Based on Student's *t* distribution (similar to normal)

$$\text{t-stat} = \frac{M(\textcolor{teal}{A}, \textcolor{blue}{x}, \textcolor{blue}{y}) - M(\textcolor{blue}{B}, \textcolor{blue}{x}, \textcolor{blue}{y})}{SE/\sqrt{m}}$$

- $m$ : size of the paired sample ( $\textcolor{blue}{x}, \textcolor{blue}{y}$ )
- $SE$ : standard deviation of difference  $d = [\hat{y}_A = \textcolor{blue}{y}] - [\hat{y}_B = \textcolor{blue}{y}]$
- P-value: look up **Student's *t* table**,  $m - 1$  degrees of freedom
  - E.g. <https://homepage.divms.uiowa.edu/~mbognar/applets/t.html>
  - In practice: `scipy stats.ttest_rel`

## T-test toy example (gory details)

$$M(\textcolor{teal}{A}, \textcolor{blue}{x}, \textcolor{blue}{y}) = \frac{1+1+1+0+1}{5} = 0.8 \quad M(\textcolor{blue}{B}, \textcolor{blue}{x}, \textcolor{blue}{y}) = \frac{1+1+0+1+0}{5} = 0.6$$

$$d = [\hat{y}_A = \textcolor{blue}{y}] - [\hat{y}_B = \textcolor{blue}{y}] = (0, 0, 1, -1, 1) \quad \bar{d} = 0.2$$

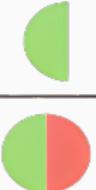
$$SE = \sqrt{\frac{(0-.2)^2 + (0-.2)^2 + (1-.2)^2 + (-1-.2)^2 + (1-.2)^2}{5-1}} = \sqrt{\frac{2.8}{4}} = 0.837$$

$$\begin{aligned}\text{t-stat} &= \frac{M(\textcolor{teal}{A}, \textcolor{blue}{x}, \textcolor{blue}{y}) - M(\textcolor{blue}{B}, \textcolor{blue}{x}, \textcolor{blue}{y})}{SE/\sqrt{m}} \\ &= \frac{0.8 - 0.6}{0.837/\sqrt{5}} \\ &= \mathbf{0.5345}\end{aligned}$$

$$\text{p-value} = P(X > 0.5345) = \mathbf{0.3107}$$

System **A** not significantly better than **B** at  $\alpha = 0.05$  level!

# Non parametric tests

$$\text{Precision} = \frac{\text{green}}{\text{red} + \text{green}}$$


$$\text{Recall} = \frac{\text{green}}{\text{green} + \text{yellow}}$$


- Problem of  $t$ -test : assumes  $M(A, x, y) \sim$  normally distributed
- Other metrics :
  - Recall  $R = \frac{tp}{tp+fn}$  ,  $tp + fn$  constant  
→  $t$ -test OK ✓
  - Precision  $P = \frac{tp}{tp+fp}$  depends on  $tp + fp$ , unknown distribution  
→  $t$ -test not OK ✗
  - F-score =  $\frac{2 \cdot P \cdot R}{P+R}$  depends on  $P$ , unknown distribution  
→  $t$ -test not OK ✗

# Parametric vs. non parametric

---

Many authors use the terms **parametric vs. non parametric tests**

- What does it mean?

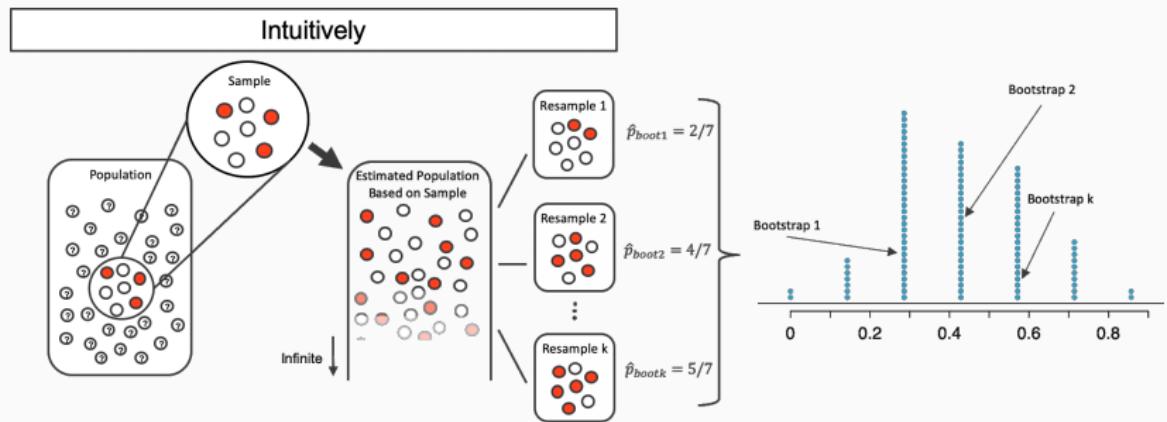
Most of the time, “parametric” means:

*The underlying random variable is **normally distributed***

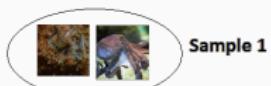
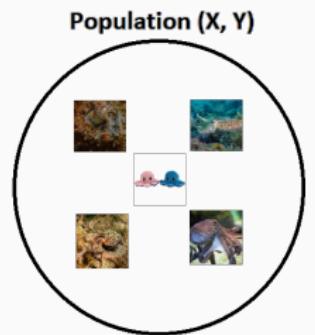
- Alternative : non parametric tests
  - 1. No sampling
    - Fast
    - Conservative,  $A$  no better than  $B$  for small  $\delta$  (low power)
    - E.g. sign test, McNemar, Wilcoxon
  - 2. With sampling
    - Slow
    - Powerful, rejects  $H_0$  more often (low type-II error probability)
    - E.g. permutation (randomised approximation), bootstrap

# Bootstrap

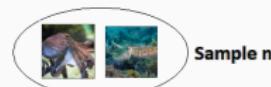
Idea : estimate  $M$  distribution by random re-sampling in  $(x, y)$



# Poolstrap



:



$$Acc(A, \text{Sample 1}) = \frac{2}{2}$$

$$Acc(B, \text{Sample 1}) = \frac{1}{2}$$

$$Acc(A, \text{Sample 2}) = \frac{1}{2}$$

$$Acc(B, \text{Sample 2}) = \frac{1}{2}$$

$$Acc(A, \text{Sample k}) = \frac{2}{2}$$

$$Acc(B, \text{Sample k}) = \frac{2}{2}$$

# Bootstrap for significance

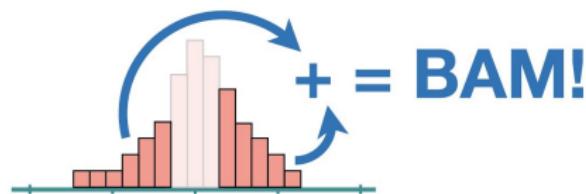
---

```
1 deltaobs = M(A,x,y) - M(B,x,y)    # delta on test set
2 R = 10000                           # 10k random samples
3 for i = 1 .. R :
4     xs, ys = sample(x,y,m)          # with repetition
5     deltasample = M(A,xs,ys) - M(B,xs,ys)
6     if deltasample > 2 * deltaobs :
7         r = r + 1
8 pvalue = r/R
```

# Why comparing with $2 \times \text{deltaobs}$ ?

- Overall sample mean is  $\delta_{A-B}(x, y)$ 
  - Standardize bootstrap sample  $(x', y')$  to assume  $H_0$  (zero effect)
  - $H_0 \implies$  subtract observed delta from each bootstrap sample value
  - $\delta_{A-B}(x', y') - \delta_{A-B}(x, y) > \delta_{A-B}(x, y)$
- Details: <https://www.youtube.com/watch?v=N4ZQQqyIf6k>

## Using Bootstrapping...

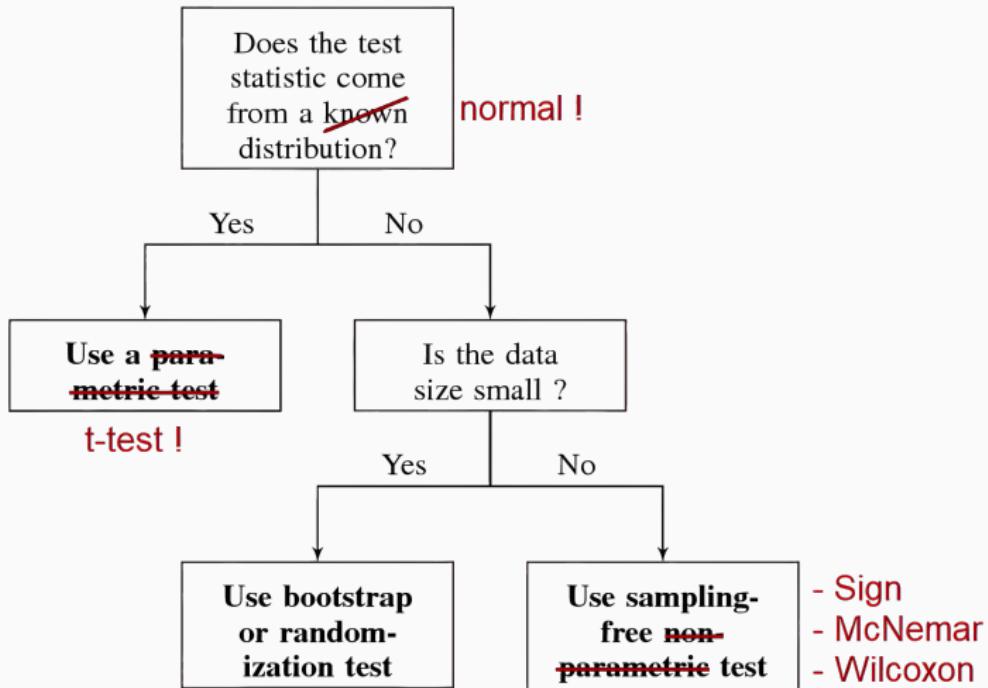


...to calculate  $p$ -values!!!

# Bootstrap and hypothesis testing

- Parametric tests were more popular in the past
  - No computer power for intensive resampling
- Today, non-parametric tests like bootstrap are most popular
  - No need to assume normality
- Bootstrap: easy system A and B comparison
  - Same test set  $(x, y)$   $\Rightarrow$  paired sample

# Which test to apply?



Source: adapted from Dror et al. (2018)

# Community's practice

## Computational linguistics conferences (ACL) and journals (TACL)

General Statistics	ACL '17	TACL '17
Total number of papers	196	37
# papers that <b>do not</b> report significance	117	15
# papers that report significance	63	18
# papers that report significance but use the <b>wrong</b> statistical test	6	0
# papers that report significance but do not mention the test name	21	3

Source: Dror et al. (2018)

# Outline

---

Statistics in a nutshell

Evaluation metrics

Statistical significance

Discussion

# Error analysis

---

- Characterise the errors in our system's output
- Scripts to print **characteristics of errors**
  - Frequency, length, resolution, predicted/gold class, ...
  - Example: elements whose prediction is furthest from gold value
- Manual error annotation: taxonomies, guidelines
  - Gain insight on most promising improvements

# Interpretability analysis

Try to understand **why** systems generate a prediction

- Feature-based methods (SHAP, LIME)
  - Which parts of the inputs influence prediction?
- Visualisation
  - Attention salience, 2-D projection (UMAP, t-SNE, topology)
- Adversarial examples, perturbations
  - Difficult minimal pairs



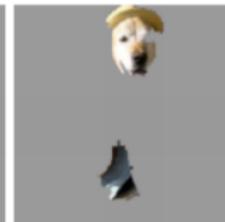
(a) Original Image



(b) Explaining *Electric guitar*



(c) Explaining *Acoustic guitar*



(d) Explaining *Labrador*

Source: <https://homes.cs.washington.edu/~marcotcr/blog/lime/>

# Negative results

- Well designed hypotheses → interesting “negative” results
- Experiments require **persistence** and some **faith**
- Source of **frustration**: publish or perish
  - Is it a problem with my results or with the system?
- Negative results are publishable if **sound experimental design**



# Leaderboards

- Remember Goodhart's law (metric  $\neq$  objective)
- Beating state of the art is good
- Learning **something interesting about the problem** is better
- From time to time: remember the research question



# Thanks!

That's all for today

---

Carlos Ramisch  
[first.last@lis-lab.fr](mailto:first.last@lis-lab.fr)

M2 IAAA - based on the course *Zen Research*  
By Carlos Ramisch and Manon Scholivet

# Sources i

- Adeline Paiement's course *Initiation à la recherche*
- Bruce et al. (2020) *Practical Statistics for Data Scientists*, 2nd Ed.
- Lanhenke (2022) *Understanding Random Variables and Probability Distributions*
- Yibi Huang's statistics course <https://www.stat.uchicago.edu/~yibi/teaching/stat220/17aut/Lectures/L22.pdf>
- Wikipedia *Precision and recall* and *Spearman Correlation*
- Youtube channels: *DATAtab*, *StatQuest*
- Feedback from participants of previous course editions
- Slides illustrated with the help of: Google images, [imgupscaler.com](http://imgupscaler.com), Canva

## Sources ii

- Slides written with the help of: ChatGPT, Google Bard, DeepL, Linguee, Overleaf
- Funding: French ANR, through SELEXINI project (ANR-21-CE23-0033-01)

Backup slides

# Random variables: formal definition i

- **Experiment:** flip 3 different coins, note head (H) or tail (T)
- The **sample space**  $S$  contains all possible experiment outcomes  
→ The subsets of  $S$  are called **events**  $E_i$
- The **random variable**  $X$  denotes the number of heads (H)
  - A variable whose exact value is unknown or irrelevant
  - We know (or estimate) its **probability distribution**  $P\{X = x_i\}$

$E_i$	$\{HHH\}$	$\{THH, HTH, HHT\}$	$\{TTH, THT, HTT\}$	$\{TTT\}$
$P(E_i)$	$1/8$	$1/8 + 1/8 + 1/8$	$1/8 + 1/8 + 1/8$	$1/8$
$X$	0	1	2	3
$P\{X = x_i\}$	$1/8$	$3/8$	$3/8$	$1/8$

## Formalisation

A **random variable** is a function  $X : S \rightarrow \mathbb{R}$  such that :

### 1. Discrete random variable:

- Its set of possible values  $X(S) = \{x_i, i \in \mathbb{N}^*\}$  is countable
- For all  $x_i \in X(S)$ :  $\{X = x_i\} \Leftrightarrow \{e_i \in S | X(e_i) = x_i\} \in \mathcal{F}$
- $\mathcal{F}$  is the set of all possible events (subsets) of  $S$
- $p(x_i) = P\{X = x_i\}$  is the **probability mass function** of  $X$

### 2. Continuous random variable:

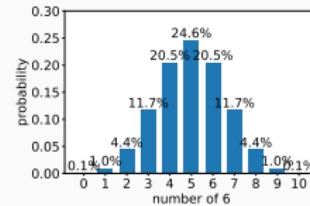
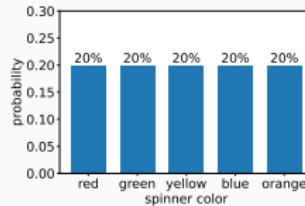
- $\forall$  value  $x \in (-\infty, +\infty)$ ,  $\forall$  interval  $B \in \mathbb{R}$
- A non-negative function  $P\{X \in B\} = \int_B f(x) dx$  exists
- $f(x)$  is the **probability density function** of  $X$

# Types of probability distributions

- Discrete random variables

→ Bar graphic, finite set of values

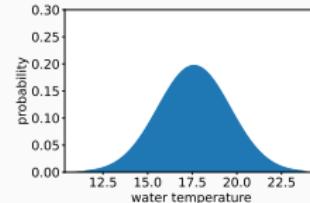
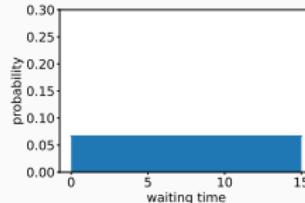
→ Probability at exact value  $P\{X = a\}$



- Continuous random variables

→ Line graphic, uncountable set of values (real numbers)

→ Probability of interval  $P\{a < X < b\}$



# Random sample or i.i.d. variables?

- Sampled items can be seen as  $n$  random variables  $X_1 \dots X_n$ 
  - For instance, tossing a coin  $n$  times
- We assume that all variables have the **same distribution**
- We assume that all items are **independent**<sup>4</sup>
- This is often stated as **independent and identically distributed**
  - The acronym **i.i.d.** is usually employed in probability

---

<sup>4</sup>Formally:  $\forall X_i \neq X_j, \forall a, b \in X_i(S) \quad P\{X_i = a | X_j = b\} = P\{X_i = a\}$

# Random sample or i.i.d. variables?

- Sampled items can be seen as  $n$  random variables  $X_1 \dots X_n$ 
  - For instance, tossing a coin  $n$  times
- We assume that all variables have the **same distribution**
- We assume that all items are **independent**<sup>4</sup>
- This is often stated as **independent and identically distributed**
  - The acronym **i.i.d.** is usually employed in probability

Random sample = set of  $n$  **values** of i.i.d. variables  $X_1 \dots X_n$

---

<sup>4</sup>Formally:  $\forall X_i \neq X_j, \forall a, b \in X_i(S) \quad P\{X_i = a | X_j = b\} = P\{X_i = a\}$

## Getting to the variance

Idea 1: average the difference between each value and the mean

$$\sum_{i=1}^n \frac{x_i - \bar{x}}{n}$$

- Calculate this amount for the sample [-4, -4, 4, 4]

## Getting to the variance

Idea 1: average the difference between each value and the mean

$$\sum_{i=1}^n \frac{x_i - \bar{x}}{n}$$

- Calculate this amount for the sample  $[-4, -4, 4, 4]$

$$\frac{(-4 - 0) + (-4 - 0) + (4 - 0) + (4 - 0)}{4} = 0 \quad \text{@}$$

## Getting to the variance

Idea 2: average the **absolute value** of the  $x_i - \bar{x}$  difference

$$\sum_{i=1}^n \frac{|x_i - \bar{x}|}{n}$$

- Calculate this amount for the sample [-4, -4, 4, 4]

## Getting to the variance

Idea 2: average the absolute value of the  $x_i - \bar{x}$  difference

$$\sum_{i=1}^n \frac{|x_i - \bar{x}|}{n}$$

- Calculate this amount for the sample  $[-4, -4, 4, 4]$

$$\frac{|-4 - 0| + |-4 - 0| + |4 - 0| + |4 - 0|}{4} = 4 \quad \text{☺}$$

## Getting to the variance

Idea 2: average the **absolute value** of the  $x_i - \bar{x}$  difference

$$\sum_{i=1}^n \frac{|x_i - \bar{x}|}{n}$$

- Calculate this amount for the sample [-6, -2, 1, 7]

## Getting to the variance

Idea 2: average the absolute value of the  $x_i - \bar{x}$  difference

$$\sum_{i=1}^n \frac{|x_i - \bar{x}|}{n}$$

- Calculate this amount for the sample [-6, -2, 1, 7]

$$\frac{|-6 - 0| + |-2 - 0| + |1 - 0| + |7 - 0|}{4} = 4 \quad \text{⊗}$$

Moreover, absolute value is not differentiable at 0

This is inconvenient: <https://www.youtube.com/watch?v=sHRBg6BhKjI>

## Getting to the variance

Idea 3: average the squared differences  $x_i - \bar{x}$

$$\sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n}$$

- Calculate this amount for the sample [-4, -4, 4, 4]

## Getting to the variance

Idea 3: average the squared differences  $x_i - \bar{x}$

$$\sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n}$$

- Calculate this amount for the sample [-4, -4, 4, 4]

$$\frac{(-4 - 0)^2 + (-4 - 0)^2 + (4 - 0)^2 + (4 - 0)^2}{4} = 64 \quad \odot$$

## Getting to the variance

Idea 3: average the squared differences  $x_i - \bar{x}$

$$\sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n}$$

- Calculate this amount for the sample [-6, -2, 1, 7]

## Getting to the variance

Idea 3: average the squared differences  $x_i - \bar{x}$

$$\sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n}$$

- Calculate this amount for the sample [-6, -2, 1, 7]

$$\frac{(-6 - 0)^2 + (-2 - 0)^2 + (1 - 0)^2 + (7 - 0)^2}{4} = 90 \quad \text{@}$$

Source: Example adapted from

<https://www.mathsisfun.com/data/standard-deviation.html>

# The larger the better

- The **expected value** of a (discrete) random variable:

$$E[X] = p(x_1)x_1 + p(x_2)x_2 + \dots + p(x_n)x_n$$

- **Sample mean**  $\bar{x} \rightarrow$  normalised sum of  $n$  i.i.d. random variables

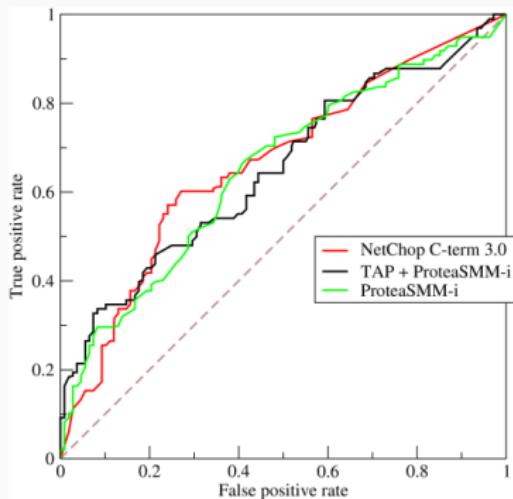
$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

- The **law of large numbers** states that  $\bar{x} \rightarrow E[X]$  for large  $n$   
→ The (sample) mean  $\bar{x}$  is an **estimator** of the expected value  $E[X]$

The larger the sample, the better  $\bar{x}$  approximates “true” mean  $E[X]$

# ROC curve

ROC curves (*Receiver Operating Characteristic*) are very useful to choose a threshold.



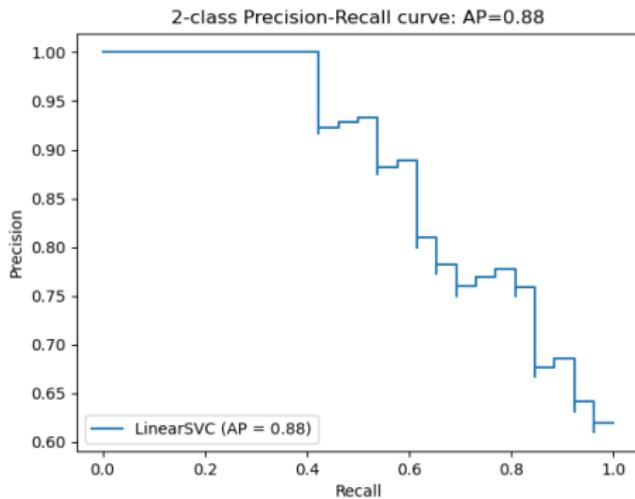
The AUC (*Area Under ROC*) is often used to estimate the model skill.

Image from Wikipedia



## Precision-recall curve

Another way to do this is to use the Precision and the Recall instead of using the True positive and the False positive rates.



# Mean average precision

- Model predicts a numerical score
- Gold class is binary or discrete
- Evaluate without setting a fixed threshold

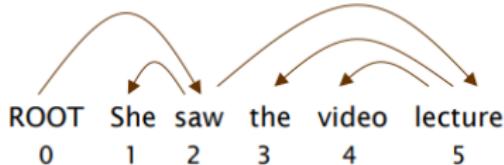
Predicted Rank	1	2	3	4	5	30	
Target	True	False	True	False	True	False	
Day	Day 1 ✓	Day 8 ✗	Day 27 ✓	Day 11 ✗	Day 29 ✓	...	Day 30 ✗
Precision	@1	@2	@3	@4	@5	@30	
=	1/1	0/2	2/3	0/4	3/5	0/30	

$AP@5 = 1/3(1/1 + 0/2 + 2/3 + 0/4 + 3/5) = 0.76$

# Structured prediction

- How to compare structured objects?
  - Sub-sequences
  - Clusters
  - Syntax trees
  - Graphs

# Structured prediction example: LAS/UAS



$$\text{Acc} = \frac{\# \text{ correct deps}}{\# \text{ of deps}}$$

"unlabelled attachment score"

$$\text{UAS} = 4 / 5 = 80\%$$

$$\text{LAS} = 2 / 5 = 40\%$$

"labelled AS"

## Gold

1	2	She	nsubj
2	0	saw	root
3	5	the	det
4	5	video	nn
5	2	lecture	obj

## Parsed

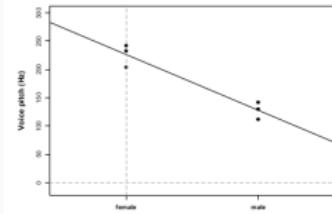
1	2	She	nsubj
2	0	saw	root
3	4	the	det
4	5	video	nsubj
5	2	lecture	ccomp

Source: <https://x-wei.github.io/xcs224n-lecture5.html>

- **Entropy:** alternative view of variability/skewness
  - $H = -\sum p(x_i) \log p(x_i)$  → amount of uncertainty
  - $H = \max$  for uniform distribution (unpredictable)
  - $H = 0$  for highly skewed distribution (predictable)
- Other useful notions:
  - Cross entropy
  - Mutual information
  - Kullbak-Leibler divergence (asymmetric)
  - Jensen–Shannon divergence (symmetric)

# Models for categorical variables

- **ANOVA:** Generalise t-test for more than 2 means
- **Linear model:** predict a linear regression slope
  - Is the slope significantly different from zero ?
  - Notation: pitch  $\approx$  sex + $\varepsilon$
- **Mixed model:** more sophisticated for multiple factors



Source: <https://bodo-winter.net/tutorials.html>

# Statistics libraries

- Visual: Excel, Libreoffice, ...
- Python: `matplotlib`, `numpy`, `scipy`, `sklearn`, ...
- R: multiple libraries including linear models
- Proprietary: Matlab, SPSS, ...